International Journal of Computer Vision 68(1), 43–52, 2006 © 2006 Springer Science + Business Media, LLC. Manufactured in The Netherlands. DOI: 10.1007/s11263-005-4841-0

Wide Baseline Matching between Unsynchronized Video Sequences

LIOR WOLF

Massachusetts Institute of Technology, The Center for Biological and Computational Learning, Cambridge, MA liorwolf@mit.edu

ASSAF ZOMET

Department of Computer Science, Columbia University, New York, NY zomet@cs.columbia.edu

Received December, 2002; Accepted February, 2005

First online version published in March, 2006

Abstract. 3D reconstruction of a dynamic scene from features in two cameras usually requires synchronization and correspondences between the cameras. These may be hard to achieve due to occlusions, different orientation, different scales, etc. In this work we present an algorithm for reconstructing a dynamic scene from sequences acquired by two uncalibrated non-synchronized fixed affine cameras.

It is assumed that (possibly) different points are tracked in the two sequences. The only constraint relating the two cameras is that every 3D point tracked in one sequence can be described as a linear combination of some of the 3D points tracked in the other sequence. Such constraint is useful, for example, for articulated objects. We may track some points on an arm in the first sequence, and some other points on the same arm in the second sequence. On the other extreme, this model can be used for generally moving points tracked in both sequences without knowing the correct permutation. In between, this model can cover non-rigid bodies with local rigidity constraints.

We present linear algorithms for synchronizing the two sequences and reconstructing the 3D points tracked in both views. Outlier points are automatically detected and discarded. The algorithm can handle both 3D objects and planar objects in a unified framework, therefore avoiding numerical problems existing in other methods.

Keywords: structure from motion, video synchronization, wide base-line matching

1. Introduction

Many traditional algorithms for reconstructing a 3D scene from two or more cameras require correspondences between the images. This becomes challenging in some cases, for example when the cameras have different zoom factors, or large vergence (wide-baseline stereo) (Tuytelaars and Van Gool, 2000; Schaffalitzky

and Zisserman, 2001). Using a moving video camera rather than a set of static cameras helps in overcoming some of the correspondence problems, but may decrease the stability and accuracy of the reconstruction. Moreover, the reconstruction from a moving camera becomes harder if not impossible when the scene is not rigid.

In this paper we present an algorithm for reconstructing a non-rigid scene from two fixed video cameras. It is assumed that feature points are tracked for each of the cameras, but no correspondences are available between the cameras. Moreover, The points

This work was done while the authors were PhD students in the School of Computer Science and Engineering, the Hebrew University of Jerusalem.

tracked by the two cameras may be different. Instead we use a weaker assumption: Every 3D point tracked in the second sequence can be described as a linear combination of some of the 3D points tracked in the first sequence. The coefficients of this linear combination are unknown but fixed throughout the sequence. Since the cameras view the same scene, this assumption is reasonable. For example, when a point tracked in the second camera belongs to some rigid part of the scene, it can be expressed as a linear combination of some other points on the same part, tracked in the first camera.

Our non-rigidity concept is quite strong. We put no *explicit* limitation on the motion of the points from one time to another. The linear combination with fixed coefficients assumption suggests that there exist local rigid structure or that some of the points, which are allowed to move freely, are tracked in both sequences (although we do not know the correspondence between the sequences).

This paper is divided to two main sections. First, we show how to synchronize two video sequences. This can be used for any application requiring camera synchronization, as well as for action indexing. Then, based on this synchronization, we show how to use the measurement from all frames simultaneously for reconstructing the non-rigid scene at different times. We conclude by experiments.

1.1. Previous Work

Using several sequences without correspondences between them was done in Caspi and Irani (2001) for alignment in space and time, and was followed by Wolf and Zomet (2002) for synchronization and selfcalibration of cameras on a stereo rig. These algorithms compute the motion of each camera independently, which is problematic for dynamic scenes. Then they combine the motions of the cameras to express the inter-camera rigidity. In Dornaika and Chung (1999), motion correspondences were used to extract stereo correspondences. It was assumed that the scene was rigid and that the geometry of the two perspective cameras is known.

A second class of related work included algorithms for reconstruction of non-rigid scenes from feature points. A solution for 3D reconstruction of k moving bodies from a sequence of affine views was presented in Costeira and Kanade (1998) and later on in Kanatani (2001) using factorization. Improved solutions were suggested in Zelnik-Manor and Irani (2003) and Vidal and Hartley (2004). These algorithms exploit the fact that there are more measurements arising from each one of the objects than the minimal number required to span this motion (usually the number of bodies is small). This redundancy is used to identify the motions. In this work the rigidity relations between points is weaker, allowing many different motions of small bodies with different dimensions. For example, local rigidity constraints of lines and planes can be exploited to express non-rigid bodies, e.g. trees and clothes. In the experiments we compared our method to Costeira and Kanade (1998), and showed that even in a k body setting as well it is worthwhile using a second camera.

Class-based methods also address reconstruction of dynamic scenes (Vetter and Poggio, 1997; Leventon and Freeman, 1998; Bregler et al., 2000, Brand, 2001). These assume that the points in 3D can be expressed as a linear combination of a small morph basis. The basis can be either provided as a learned prior (Leventon and Freeman, 1998; Vetter and Poggio, 1997), or computed by the algorithm (Bregler et al., 2000; Brand, 2001). In this work, rather than expressing all the points as linear combinations of some morph basis, we set some of the points as linear combination of other points, thus expressing local rigidity of subsets of the features. Still, the assumptions of the class-based approach and the presented work can be expressed in a similar way: The matrix containing the positions of the 3D points in different times has a low rank. Therefore, a single-camera factorization algorithm designed for class based setting (Bregler et al., 2000) can be adapted also for reconstruction in our proposed setting.

The method proposed in this paper uses two cameras and can be seen as a compromise between having the prior 3D model given to the algorithm (Vetter and Poggio, 1997; Leventon and Freeman, 1998) and computing the model from the data using factorization (Bregler et al., 2000; Brand, 2001). By adding a second camera we obtain two main advantages over the factorization methods. First, in factorization methods there is an ambiguity expressed in an unknown invertible matrix between the two factors. In contrast, in our approach, the ambiguity is resolved. Second, the proposed method is robust to outliers as it includes a way to discard outliers by testing each feature separately. Using factorization on data with outliers degrades the quality of the results. We elaborate more on this point in Section 4.1.

2. Formal Statement

Assume two affine cameras $(2 \times 4 \text{ matrices})$ view different 3D points across time $1 \le t \le T$. Let $P_j^{(t)} \in \Re^3$, $j = 1 \dots n$, be the 3D points tracked by the first camera.

$$\begin{pmatrix} x_i^{(t)} \\ y_i^{(t)} \end{pmatrix} = \begin{pmatrix} a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \end{pmatrix} P_i^{(t)} + \begin{pmatrix} a_4 \\ b_4 \end{pmatrix} \quad (1)$$

Instead of looking directly on the point measurements $(x, y)^{\top}$ we look on the motion relative to some reference frame:

$$\begin{pmatrix} u_i^{(t)} \\ v_i^{(t)} \end{pmatrix} = \begin{pmatrix} a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \end{pmatrix} \begin{pmatrix} P_i^{(t)} - P_i^{(0)} \end{pmatrix}$$
$$= \begin{pmatrix} a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \end{pmatrix} dP_i^{(t)},$$
(2)

where $dP_i^{(t)} = P_i^{(t)} - P_i^{(o)}$ is the motion of point *i* from the time of the reference frame (the first frame) to the time of frame *t*, and $u_i^{(t)}$, $v_i^{(t)}$ denote the motion of point *i* in the image.

Here Q_i is the vector specifying the linear combination coefficients, such that $\sum_i Q_i(j) = 1$.

Note that this linear combination applies also to the 3D displacements $d\hat{P}_i^{(t)}$ and so:

$$\begin{pmatrix} \hat{u}_{i}^{(t)} \\ \hat{v}_{i}^{(t)} \end{pmatrix} = \begin{pmatrix} \hat{a}_{1} & \hat{a}_{2} & \hat{a}_{3} \\ \hat{b}_{1} & \hat{b}_{2} & \hat{b}_{3} \end{pmatrix} \times \begin{pmatrix} dP_{1}^{(t)} & dP_{2}^{(t)} & \dots & dP_{n}^{(t)} \end{pmatrix} Q_{i}$$
(5)

$$= \begin{pmatrix} \hat{a}_{1} & \hat{a}_{2} & \hat{a}_{3} \\ \hat{b}_{1} & \hat{b}_{2} & \hat{b}_{3} \end{pmatrix} \begin{pmatrix} u_{1}^{(t)} & u_{2}^{(t)} & \dots & u_{n}^{(t)} \\ v_{1}^{(t)} & v_{2}^{(t)} & \dots & v_{n}^{(t)} \\ dZ_{1}^{(t)} & dZ_{2}^{(t)} & \dots & dZ_{n}^{(t)} \end{pmatrix} Q_{i}$$
(6)

Where $dZ_i^{(t)}$ are the displacements along the Z axis of the points tracked in the first sequence and the last equality is by substituting the projection equation of the first camera.

Reorganizing the terms we get:

$$\begin{pmatrix} \hat{u}_{i}^{(t)} \\ \hat{v}_{i}^{(t)} \end{pmatrix} = \begin{pmatrix} Q_{i}(1)\hat{a}_{1} & Q_{i}(1)\hat{a}_{2} & Q_{i}(1)\hat{a}_{3} & Q_{i}(1)\hat{a}_{3} \\ Q_{i}(1)\hat{b}_{1} & Q_{i}(1)\hat{b}_{2} & Q_{i}(1)\hat{b}_{3} & Q_{i}(1)\hat{b}_{3} \end{pmatrix}$$

We will use the coordinate system of the first camera as our world (3D) coordinate system and so our projections onto the two cameras are given by:

$$\begin{pmatrix} u_i^{(t)} \\ v_i^{(t)} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix} dP_i^{(t)}, \begin{pmatrix} \hat{u}_i^{(t)} \\ \hat{v}_i^{(t)} \end{pmatrix} = \begin{pmatrix} \hat{a}_1 & \hat{a}_2 & \hat{a}_3 \\ \hat{b}_1 & \hat{b}_2 & \hat{b}_3 \end{pmatrix} d\hat{P}_i^{(t)},$$
(3)

where we denote points related to the second camera using a similar notation with a hat.

The assumption we make in this work is that every 3D point tracked in the second sequence $\hat{P}_i^{(t)}$, $1 \ge i \ge m$ can be expressed as a linear combination of the points tracked in the first sequence:

$$\hat{P}_{i}^{(t)} = \begin{pmatrix} P_{1}^{(t)} & P_{2}^{(t)} & \dots & P_{n}^{(t)} \end{pmatrix} Q_{i}.$$
 (4)

$$\begin{array}{ccccc} {}_{i}(2)\hat{a}_{1} & Q_{i}(2)\hat{a}_{2} & \dots & Q_{i}(n)\hat{a}_{3} \\ {}_{i}(2)\hat{b}_{1} & Q_{i}(2)\hat{b}_{2} & \dots & Q_{i}(n)\hat{b}_{3} \end{array} \begin{pmatrix} u_{1}^{(t)} \\ v_{1}^{(t)} \\ Z_{1}^{(t)} \\ \vdots \\ u_{n}^{(t)} \\ v_{n}^{(t)} \\ Z_{n}^{(t)} \end{pmatrix}$$
(7)

Let $A_i = (Q_i(1)\hat{a}_1, Q_i(1)\hat{a}_2, Q_i(1)\hat{a}_3, Q_i(2)\hat{a}_1, Q_i(2)\hat{a}_2, \ldots, Q_i(n)\hat{a}_3)$, and let B_i be the similar expression involving $\hat{b}_1, \hat{b}_2, \hat{b}_3$. Let \hat{M} be the matrix whose rows are $A_1, \ldots, A_m, B_1, \ldots, B_m$. and let *C* be the matrix whose columns are $(u_1^{(t)}, v_1^{(t)}, Z_1^{(t)}, \ldots, u_n^{(t)}, v_n^{(t)}, Z_n^{(t)})$ for $1 \le t \le T$.

Note that \hat{m} is point-dependent, but does not vary with time, whereas the matrix *C* is time-dependent, but does not depend on points in the second sequence.

Let \hat{U} and \hat{V} be the matrices containing the flow in the second sequences:

$$\hat{U} = \begin{bmatrix} \hat{u}_{1}^{(1)} & \dots & \hat{u}_{m}^{(1)} \\ \vdots & \ddots & \vdots \\ \hat{u}_{1}^{(T)} & \dots & \hat{u}_{m}^{(T)} \end{bmatrix}, \quad \hat{V} = \begin{bmatrix} \hat{v}_{1}^{(1)} & \dots & \hat{v}_{m}^{(1)} \\ \vdots & \ddots & \vdots \\ \hat{v}_{1}^{(T)} & \dots & \hat{v}_{m}^{(T)} \end{bmatrix}$$

46 Wolf and Zomet

Stacking all Eq. (7) we get:

$$\begin{pmatrix} \hat{U}^T\\ \hat{V}^T \end{pmatrix}_{2m\times T} = \hat{M}_{2m\times 3n} C_{3n\times T}$$
(9)

A similar expression can be derived for the first camera as well:

$$\begin{pmatrix} U^T \\ V^T \end{pmatrix}_{2n \times T} = M_{2n \times 3n} C_{3n \times T}$$
 (10)

These equations will be used in both the synchronization algorithm (Section 3) and the 3D reconstruction algorithm (Section 4).

3. Synchronization and Action Indexing

Given two sequences of a dynamic scene in the above setting, they are first synchronized. Having synchronization, we later describe how to compute the structure of the scene across time. The ability to synchronize two sequences of the same action can be used for other dynamic-scene applications such as action indexing. Two matching actions should have a minimal synchronization error.

Consider the matrix obtained by stacking the flow in both sequences side by side: $E = (U, V, \hat{U}, \hat{V})$. This is a matrix of size $T \times (2n + 2m)$ and since we consider long sequences we expect its rank to be (2n + 2m). However, this holds only when there is no connection between the sequences, i.e when they are not synchronized. Combining Eqs. (9) and (10) we get:

$$E^{T} = \begin{pmatrix} M \\ \hat{M} \end{pmatrix}_{(2m+2m)\times 3n} C_{3n\times T}$$
(11)

therefore the rank of E is bounded by 3n. We propose to synchronize the sequences by examining the rank of the matrix E for various temporal shifts between the sequences.

The upper bound in Eq. (11) is usually not tight, depending on the rigidity of the scene. For example, when the scene contains k rigid bodies the rank is bounded by 3k. Still, we expect the rank of the matrix E to decrease in the synchronized case at least as much as it decreases in the unsynchronized case.

Outliers (points which do not satisfy our assumptions even remotely) pose a different kind of problem for this rank bound, since they may increase the rank of the matrix E. However for outliers we expect the

rank of the matrix E to increase both for the synchronized case and for the unsynchronized case by the same amount. Therefore in both cases (partially rigid scenes and outliers) we expect the rank in the synchronized case to be lower than the rank in the unsynchronized case).

In practice, however, synchronizing by examining the rank of a matrix is problematic since due to noise our matrices will almost always be of full rank. In order to deal with this we propose two solutions. The first solution examines the *effective rank* of the measurements. The second solution examines the principal angles between the space of measurements of both sequences.

The first solution, instead of bounding the rank of *E* by 3*n*, uses a heuristic to define \hat{n} , the *effective rank* of (*U*,*V*). We compute the singular values s_1 ,\ldots, s_{h1} of (*U*,*V*), and set $N = \operatorname{argmin}_j \{\sum_{k=1}^{2j} s_k > J\}$ for some threshold *J* (we used $J = 0.99 \sum_{k=1}^{h_1} s_k$). This is equivalent to finding the minimal rank of a matrix in a given error bound under Frobenius norm.

In order to synchronize we select the synchronization for which the algebraic measure g(E) defined below is minimized. The minimization is carried out over all possible synchronizations where each possible synchronization provides a different matrix *E*. Let e_1, \ldots, e_{h_2} be all the singular values of *E*. Then *g* is defined by:

$$g(E) = \sum_{k=3\hat{n}+1}^{n_2} e_k$$
(12)

g measures the amount of energy in the matrix beyond the expected rank bound. It is expected that if the rank bound on E is violated by noise, this measure will be low, and otherwise, if it is broken by non-synchronized data, this measure will be high.

The second solution is based on the idea that for synchronized sequences the linear subspaces spanned by the columns of (U,V) and by the columns of (U', V') intersect. The rank of the intersection is 2n.

A stable way of measuring the extent of intersection between two linear subspaces is by using the principal angles (Golub and Loan, 1996). Let U_A , U_B represent two linear subspaces of \mathbb{R}^n . The principal angles $0 \le \theta_1 \le \cdots \le \theta_k \le (\pi/2)$ between the two subspaces are uniquely defined as:

$$cos(\theta_k) = \max_{\mathbf{u} \in U_A} \max_{\mathbf{v} \in U_B} \mathbf{u}^\top \mathbf{v}$$

subject to: $\mathbf{u}^\top \mathbf{u} = \mathbf{v}^\top \mathbf{v} = 1$, $\mathbf{u}^\top \mathbf{u}_i = 0$, $\mathbf{v}^\top \mathbf{v}_i = 0$,
 $i = 1, \dots, k - 1$

A convenient way to compute the principal angles is via the QR factorization, as follows. Let $A = Q_A R_A$ and $B = Q_B R_B$ where Q is an orthonormal basis of the respective subspace and R is a upper-diagonal $k \times k$ matrix with the Gram-Schmidt coefficients representing the columns of the original matrix in the new orthonormal basis. The singular values $\sigma_1, \ldots, \sigma_k$ of the matrix $Q_A Q_B$ are the principal angles $cos(\theta_i) = \sigma_i$.

As mentioned above the linear subspaces spanned by the columns of [U, V] and of [U', V'] intersect and the rank of the intersection is 2n. Hence, the first 2n principal angles between these column spaces vanish. In practice we are affected by noise and expect those first principal angles to be close to zero. In order to achieve synchronization we consider a function of the first few principal angles (e.g. their average). This function is computed for every possible displacement and the synchronization is chosen as to minimize this function.

3.1. Direct Synchronization using Brightness Measurements

The synchronization procedures described above are based on having a set of points tracked over time in each one of the video sequences. The accuracy of our results is therefore affected by the accuracy of the tracker. Most trackers have difficulties maintaining accurate positions for all tracked points over time, especially in dynamic scenes.

To overcome this, we present a method that uses brightness measurements instead of tracked points, following the work of Irani (1999). In Irani (1999) it was shown that the matrices U and V are closely related to matrices that are computed directly from the brightness values of the input frames I(t), $0 \le t \le T$. Specifically, they are computed from the spatial derivatives I_x and I_y of the reference image (t = 0) and the frame differences $I_t^j = I(j) - I(0)$:

$$[U,V]_{T\times 2n} \begin{bmatrix} \bar{A}, \ \bar{B} \\ \bar{B}, \ \bar{C} \end{bmatrix}_{2n\times 2n} = [G,H]_{T\times 2n}$$
(13)

Where \bar{A} , \bar{B} and \bar{C} are diagonal $n \times n$ matrices with the diagonal elements being $\bar{A}_{ii} = \sum_k (I_x(k))^2$, $\bar{B}_{ii} = \sum_k (I_x(k)I_y(k))$ and $\bar{C}_{ii} = \sum_k (I_y(k))^2$ respectively, and the summation is over a small neighborhood around pixel *i*. The elements of *G* and *H* are similarly given by $H_{ij} = -\sum_k I_x(k)I_t^i(k)$ and $H_{ij} = -\sum_k I_y(k)I_t^i(k)$. In practice we compute matrices G and H only for those pixels for which there is enough gradient information in their neighborhood, hence the matrix $\begin{bmatrix} \bar{A}, \bar{B} \\ \bar{B}, \bar{C} \end{bmatrix}$ is invertible and the column space of U and V is the same as the column space of G and H respectively. A similar constraint connects the matrices in the second sequence U' and V' to measurement matrices G' and H' computed on this sequence. We can therefore make the observations below, which are analogous to the observations we made earlier on the ranks of U, V, U', V'. In the synchronized case:

- The effective rank *N* of the matrix [*U*, *V*] is the same as the effective rank of the matrix [*G*, *H*].
- The effective rank of the matrix [G, H, G', H'] is bounded by 3N.
- The column spaces of the matrix [G, H] and of the matrix [G'.H'] intersect in a linear subspace of rank N.

Based on these observations, we use synchronization methods analogous to those presented in section 3. A care should be taken as for the number of frames used. The constraint in Eq. (13) assumes infinitesimal motion. This assumption may hold only a few frames away from the reference frame. In practice we divide both sequences to small continuous fragments (small windows in time) of length up to 15 frames. We set the reference frame as the middle frame of each fragment, and compute the measurement matrices G and H for each such fragment separately. We then compare each fragment in the first sequence to each fragment in the second sequence, and compute an algebraic error.

Let k, k' be the number of fragments we extract from each sequence. The comparison of fragments between both sequences results in a matrix of size $k \times k'$. A synchronization will appear as the line in this matrix which has the minimum average value . The offset of the line determines the shift of the synchronization and the slope of it determines the ratio of frame rates between the sequences.

4. Reconstruction

Given the synchronization of the two sequences, the next stage is 3D reconstruction. Let $D = null(\begin{bmatrix} U^T \\ V^T \end{bmatrix})$ be a matrix whose columns are all orthogonal to the rows of $\begin{bmatrix} U^T \\ V^T \end{bmatrix}$ The rows of $\begin{bmatrix} U^T \\ V^T \end{bmatrix}$ are the rows of *C* which correspond to motion along the *X* and *Y* axes (see Eq. (9)). Hence the first two out of every 3 rows in *CD*

48 Wolf and Zomet

vanish (all which is left is the information regarding the *Z* coordinates of the displacements).

By Multiplying both sides of eq. (9) with D we define a matrix K:

$$K = \begin{pmatrix} \hat{U} \\ \hat{V} \end{pmatrix} D = \hat{M}CD$$
$$= \begin{pmatrix} \hat{a}_3 Q_1 \\ \hat{b}_3 Q_1 \\ \vdots \\ \hat{a}_3 Q_m \\ \hat{b}_3 Q_m \end{pmatrix} \begin{pmatrix} Z_1^{(1)} & \dots & Z_1^{(T)} \\ \vdots & \ddots & \vdots \\ Z_n^{(1)} & \dots & Z_n^{(T)} \end{pmatrix} D \quad (14)$$

Observe that each odd row of *K* equals the next row of *K* multiplied by the ratio $r = \hat{a}_3/\hat{b}_3$. Hence this ratio can be recovered by dividing the rows. For added robustness we take the median of all the measurements of *r* from all points. In Section 4.1 we elaborate on the robustness of the algorithm.

Consider the vector $l = (1 - r)^{\top}$. *l* is a direction orthogonal to the projection of the *Z* axis to the second image. Using this direction we eliminate the unknown depths of the points.

We multiply both sides of eq. (6) with l^{\top} and get:

$$l^{\top} \begin{pmatrix} \hat{u}_{i}^{(t)} \\ \hat{v}_{i}^{(t)} \end{pmatrix} = (c_{1} \ c_{2}) \begin{pmatrix} u_{1}^{(t)} & u_{2}^{(t)} & \dots & u_{n}^{(t)} \\ v_{1}^{(t)} & v_{2}^{(t)} & \dots & v_{n}^{(t)} \end{pmatrix} Q_{i}$$
(15)

Where $c_1 = l^{\top}(\hat{a}_1, \hat{b}_1), c_2 = l^{\top}(\hat{a}_2, \hat{b}_2)$, and the third row of the projection matrix just vanished $(l^{\top}(\hat{a}_3, \hat{b}_3) = 0)$.

This is a bilinear system in the unknowns c_1 , c_2 and Q_i . For each *i* we have $\rightarrow T$ equations. We solve this system by linearization, i.e by defining a vector of unknowns $u_i = (c_1Q_i, c_2Q_i)$ and converting the equations into a linear system.

From u_i recovering Q_i up to scale is is done by factoring the elements of u_i as a $2 \times n$ rank 1 matrix.



Figure 1. Sequences synchronization for cameras with different zoom factors (Fig. 1(a) and for cameras with large vergence (Fig. 1(b). Corresponding frames from the sequences of the first graph are shown in c) and d). The first graph reflects the oscillations in the input motion (walking people). The sequences of the second graph are the same as in Fig. 4 The graphs show the algebraic error g vs. frame offset.



(d)

Sequences synchronization for cameras with large vergence using the direct method (Section 3.1). Panel (a) shows the distance Figure 2. (minimal principal angle, see text) between temporal windows at different temporal shifts in the two sequences. By averaging diagonal directions in this matrix, one achieves a score for different synchronizations, as shown in panel (b). The cyclic nature of the motion can be viewed in the graph. Corresponding frames from the sequences of the first graph are shown in (c) and (d).

 c_1, c_2 > can be recovered using all of u_i , since they are not point-dependent. Then the scale is adjusted such that $\sum_{i} Q_i = 1$.

Once the Q_i are recovered we can recover the points in the first image corresponding to the points tracked in the second image:

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \hat{P}_i^{(t)}$$

$$= \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} P_1^{(t)} & P_2^{(t)} & \dots & P_n^{(t)} \end{pmatrix}$$

$$Q_{i} = \begin{pmatrix} x_{1}^{(t)} & x_{2}^{(t)} & \dots & x_{n}^{(t)} \\ y_{1}^{(t)} & y_{2}^{(t)} & \dots & y_{n}^{(t)} \end{pmatrix} Q_{i}$$
(16)

The problem of reconstructing the points tracked in the second image $\hat{P}_i^{(t)}$ is hence reduced to a simple stereo problem in each frame. Reconstructing $P_i^{(t)}$ from $\hat{P}_i^{(t)}$ is possible by taking the pseudo-inverse of the matrix whose columns are

 $Q_i, i = 1 \dots m.$



Figure 3. Feature reprojection between two sequences of a non-rigid body.

4.1. Outlier Rejection

The basis of our algorithm and many similar feature based-algorithms is the assumption that the input points fit some model and were tracked correctly along the sequence. Outlier points, e.g. points which were not tracked properly, may have in some cases a strong effect on the accuracy of the results. In this section we show that outlier points tracked in the first camera have negligible effect on the results, and outlier points tracked the second camera can be detected automatically. Hence the algorithm is robust to outlier tracks from both sequences. The first step in the algorithm consists of computing the matrix D, i.e. finding the subspace orthogonal to the 2D trajectories of all points tracked in the first camera. In case of outlier points, this subspace may be reduced, but still the resulting subspace is orthogonal to all inlier trajectories. Therefore the computation of $r = \hat{a}_3/\hat{b}_3$ described in the previous section is not influenced by outliers in the first sequence. As for the computation of Q_i —3D points in the second camera are a combination of *some* of the inliers 3D points tracked in the first camera. The coefficients of Q_i corresponding to the outlier points vanish. Outliers tracks in the first sequence therefore have little effect as long as there are not too many of them. We show next how to eliminate outliers tracks from the second sequence. For each point in the second image, several measurements of the ratio r can be computed. A simple test, e.g. by the measuring variance of this value for each point and across the points can be used to reject outliers.

In our experiments usually $2n > T > 3\hat{n}$, i.e. the matrix \hat{B} had rank *T* due to noise, but the effective rank 2*N* was smaller. For numerical reasons we chose *D* to be the vector corresponding to the smallest singular value, and therefore every point in the second image had a single measurement of the ratio *r*. Outlier points were chosen by computing the median of the *r* values of all points, and discarding points with *r* value far from the median.

5. Experiments

We have tested our algorithm on scenes of moving people. The first set of experiments tested the camera synchronization application. Two video sequences were captured by two unsynchronized cameras viewing the same non-rigid scene. In the first experiment we used tracked points, and examined the algebraic error g, as explained in Section 3, at different time shifts, and chose the shift with the minimal error. Results are presented in Fig. 1. In the second experiment we have used the direct method described in Section 3.1. Since this method assumes small image motion, we applied it to small temporal windows w(t) (typically of 15 consecutive frames), and smoothed and decimated the images. For each such pair of windows $w(t_1), w(t_2)$ from the two sequences, we measured the smallest principal angles between the subspaces defined on these windows. A typical result of this process is shown in Fig. 2(a). The matrix contains for each pair of times $\rightarrow t_1, t_2$ the principal angles between the linear spaces computed for $w(t_1)$ on the first sequence and $w(t_2)$ on the second sequence. Since both $\rightarrow t_1, t_2$ increment by the same shift, the optimal synchronization is visualized by a 45° dark line in fig. 2(a). The mean value along 45° lines is shown in Fig. 2(b). Corresponding frames of the two synchronized sequences are shown in Fig. 2-(c-f). We have confirmed these results by manually synchronizing the input sequences.

The second experiment tested the reprojection stage. Using the proposed algorithm we established correspondences between the two cameras, by transfer-



Figure 4. Feature clustering for k-body scene, a comparison with the Costeira-Kanade algorithm. Fig. 4(a) and (b) show examples of input images taken by the two cameras at different times. Note the large vergence between the cameras. Fig. 4(c) and (d) show the two clusters found by the method in this paper. Fig. 4(e) and (f) show the two clusters found by the Costeira-Kanade K-body factorization algorithm.

ring tracked points from one sequence to the other sequence. The results are presented in Fig. 3. Unfortunately, more than 50% of the tracks were not very good, and our algorithm was not able to handle these tracks as we hoped. So while in the synthetic experiments the algorithm proved to be robust to outliers, in the presented experiment we have chosen good tracks in both sequences manually.

5.1 Clustering

In order to show the benefits of using our algorithm for the settings of independently moving rigid bodies we implemented the algorithm of Costeira and Kanade (1998) and applied it to each of the two sequences separately. Based on the results we clustered the points to different rigid objects, thus comparing the results on real noisy data. The clustering method for both algorithms was based on an affinity matrix of the points. In our algorithm an affinity matrix can

52 Wolf and Zomet

be defined by the normalized correlation between the coefficient vectors Q_i . In Costeira and Kanade (1998) an affinity matrix was defined in a similar manner. As pointed out in Kanatani (2001), once the affinity matrix is defined, the choice of the clustering criteria is arbitrary. We have used a software package (ClusteringProgram http://odur.let.rug.nl/kleiweg/clustering/ clustering.html.) including several clustering algorithms, and chose the Ward method which brought the best results for both algorithms. The results, presented in Fig. 4, show that the use of an additional camera in our algorithm improves the results.

6. Discussion

Finding correspondences is challenging, especially between views that are separated by a wide baseline. In this work we have addressed the correspondence problem in a dynamic scene. We showed that by using a video rather than a single image, it is easier to resolve the correspondence problem. This is done by tracking points in each video sequence, and by using a simple and realistic constraint to relate the points tracked in the two video sequences. The proposed algorithm is simple, linear and robust to outliers. It combines the measurements from all video frames in a single computation, thus minimizing the error in all frames simultaneously. All in all, it provides a simple and robust solution for reconstructing a dynamic 3D scene from non-synchronized video cameras.

Note

 We assume T > 2N, otherwise synchronization is not possible. If T < 2n, we take D to be the vector corresponding to the smallest singular value.

References

- Available at http://odur.let.rug.nl/\kleiweg/clustering/clustering. html.
- Brand, M., 2001. Morphable 3d models from video. In *Computer Vision and Pattern Recognition*, Kauai, Hawaii, pp. II:456–463.

- Bregler, C., Hertzmann, A., and Biermann, H. 2000. Recovering non-rigid 3d shape from image streams. In *Computer Vi*sion and Pattern Recognition, Hilton Head, SC, pp. II:690– 696.
- Caspi, Y., and Irani, M. 2001. Alignment of non-overlapping sequences. In *International Conference on Computer Vision*, Vancouver, BC, pp. II: 76–83.
- Costeira, J. and Kanade, T. 1998. A multibody factorization method for independently moving-objects. *Int. J. of Computer Vision*, 29(3):159–179.
- Dornaika, F. and Chung, R. 1999. Stereo correspondence from motion correspondence. In *Computer Vision and Pattern Recognition*, pp. I:70–75.
- Golub, G. and Loan, C.V. 1996. *Matrix Computations*, 3rd ed., Johns Hopkins University Press: Baltimore, MD.
- Irani, M. 1999. Multi-frame optical flow estimation using subspace constraints. In *IEEE International Conference on Computer Vi*sion (ICCV), Corfu.
- Kanatani, K. 2001. Motion segmentation by subspace separation and model selection. In *International Conference on Computer Vision*, Vancouver, BC, pp. II: 586–591.
- Leventon, M. and Freeman, W. 1998. Bayesian estimation of 3d human motion. technical report tr 98-06, mitsubishi electric research labs.
- Schaffalitzky, F. and Zisserman, A. 2001. Viewpoint invariant texture matching and wide baseline stereo. In *International Conference on Computer Vision*, Vancouver, BC, pp. II: 636–643.
- Sturm, P. and Triggs, B. 1996. A factorization based algorithm for multi-image projective structure and motion. In *European Conference on Computer Vision*, pp. II:709– 720.
- Tomasi, C. and Kanade, T. 1992. Shape and motion from image streams under orthography: A factorization method. *Int. J. of Computer Vision*, 9(2):137–154.
- Tuytelaars, T. and Van Gool, L. 2000. Wide baseline stereo matching based on local, affinely invariant regions. In *British Machine Vision Conference*.
- Vetter, T. and Poggio, T. 1997, Linear object classes and image synthesis from a single example image. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(7):733– 742.
- Vidal, R. and Hartley, R. 2004. Motion segmentation with missing data using power factorization and GPCA. In *Computer Vision* and Pattern Recognition.
- Zelnik-Manor, L. and Irani, M. 2003. Degeneracies, dependencies and their implications in multi-body and multisequence factorizations. In *Computer Vision and Pattern Recognition*.
- Wolf, L. and Zomet, A. 2002. Sequence to sequence self calibration In *European Conference on Computer Vision*, Copenhagen, Denmark, pp. II: 370–377.